Distributed Systems

Group Membership and View Synchronous Communication

Prof R. Guerraoui Distributed Programming Laboratory



Who is there?





 In many distributed applications, processes need to know which processes are *participating* in the computation and which are not

 Failure detectors provide such information; however, that information is *not coordinated* (see next slide) even if the failure detector is perfect





• To illustrate the concept, we focus here on a group membership abstraction to coordinate the information about *crashes*

 In general, a group membership abstraction can also typically be used to coordinate the processes *joinning* and *leaving* explicitly the set of processes (i.e., without crashes)

- Like with a failure detector, the processes are informed about failures; we say that the processes install views
- Like with a perfect failure detector, the processes have accurate knowledge about failures
- Unlike with a perfect failure detector, the information about failures are coordinated: the processes install the same sequence of views

Memb1. Local Monotonicity: If a process installs view (j,M) after installing (k,N), then j > k and M < N

Memb2. Agreement: No two processes install views (j,M) and (j,M') such that $M \neq M'$

Memb3. Completeness: If a process p crashes, then there is an integer j such that every correct process eventually installs view (j,M) such that $p \notin M$

Memb4. Accuracy: If some process installs a view (i,M) and $p \notin M$, then p has crashed

Events

Indication: <membView, V>

• Properties:

Memb1, Memb2, Memb3, Memb4

Algorithm (gmp)

Implements: groupMembership (gmp).

Vses:

- PerfectFailureDetector (P).
- UniformConsensus(Ucons).
- r upon event < Init > do
 - view := (0,S);
 - correct := S;
 - wait := true;

Algorithm (gmp – cont'd)

r upon event < crash, pi > do

correct := correct \ {pi};

- upon event (correct < view.memb) and (wait = false) do
 - wait := true;
 - trigger<ucPropose,(view.id+1,correct) >;

Algorithm (gmp – cont'd)

- upon event < ucDecided, (id, memb)> do
 - view := (id, memb);
 - wait := false;
 - trigger < membView, view>;

Algorithm (gmp)



Group Membership and Broadcast



14

 View synchronous broadcast is an abstraction that results from the combination of group membership and reliable broadcast

• *View synchronous broadcast* ensures that the delivery of messages is coordinated with the installation of views

Besides the properties of *group membership* (*Memb1-Memb4*) and *reliable broadcast* (*RB1-RB4*), the following property is ensured:

VS: A message is vsDelivered in the view where it
is vsBroadcast

Events

Request:

vsBroadcast, m>

- Indication:
 - <vsDeliver, src, m>
 - <vsView, V>

If the application keeps *vsBroadcasting* messages, the *view synchrony* abstraction might never be able to *vsInstall* a new view; the abstraction would be impossible to implement

We introduce a specific event for the abstraction to **block** the application from **vsBroadcasting** messages; this only happens when a process crashes

- *Events*
 - Request:
 - vsBroadcast, m>; <vsBlock, ok>
 - Indication:
 - vsDeliver, src, m>; <vsView, V>; <vsBlock>

Algorithm (vsc)

Implements: ViewSynchrony (vs).

C Uses:

- GroupMembership (gmp).
- TerminatingReliableBroadcast(trb).
- BestEffortBroadcast(beb).

r upon event < Init > do

- \checkmark view := (0,S); nextView := \bot ;
- \checkmark pending := delivered := trbDone := \varnothing ;
- flushing := blocked := false;

- upon event <vsBroadcast,m> and (blocked = false)
 do
 - \checkmark delivered := delivered $\cup \{ m \};$
 - r trigger <vsDeliver, self, m>;
 - r trigger <bebBroadcast, [Data,view.id,m>;

upon event < bebDeliver,src,[Data,vid,m]) do
</pre>

- ✓ If(view.id = vid) and (m ∉ delivered) and (blocked = false) then
 - \checkmark delivered := delivered \cup \langle m \rangle
- r trigger <vsDeliver, src, m >;

- r upon event < membView, V > do
 - addtoTail (pending, V);
- ✓ upon (pending $\neq \emptyset$) and (flushing = false) do
 ✓ nextView := removeFromhead (pending);
 - flushing := true;
 - rtrigger <vsBlock>;

- Upon <vsBlockOk> do
 - blocked := true;
 - \checkmark trbDone := \varnothing ;
 - rtrigger <trbBroadcast, self, (view.id,delivered)>;

- Upon <trbDeliver, p, (vid, del)> do
 - *r* trbDone := trbDone $\cup \{p\}$;
 - ✓ forall m ∈ del and m ∉ delivered do
 ✓ delivered := delivered $\cup \{m\}$;
 - f trigger <vsDeliver, src, m >;

- r Upon (trbDone = view.memb) and (blocked = true) do
 - view := nextView;
 - flushing := blocked := false;
 - \checkmark delivered := \varnothing ;
 - r trigger <vsView, view>;

Consensus-Based View Synchrony

Instead of launching parallel instances of TRBs, plus a group membership, we use one consensus instance and parallel broadcasts for every view change

Roughly, the processes exchange the messages they have delivered when they detect a failure, and use consensus to agree on the membership and the message set

Algorithm 2 (vsc)

Implements: ViewSynchrony (vs).

C Uses:

- UniformConsensus (uc).
- BestEffortBroadcast(beb).
- PerfectFailureDetector(P).

- r upon event < Init > do
 - r view := (0,S);
 - r correct := S;
 - flushing := blocked := false;
 - \checkmark delivered := dset := \varnothing ;

- r upon event <vsBroadcast,m) and (blocked = false) do</pre>
 - \checkmark delivered := delivered $\cup \{ m \}$
 - r trigger <vsDeliver, self,m>;
 - r trigger <bebBroadcast,[Data,view.id,m] >;

- upon event < bebDeliver,src,[Data,vid,m]) do
 </pre>
 - f (view.id = vid) and (m ∉ delivered) and (blocked
 = false) then
 - \checkmark delivered := delivered \cup \langle m \rangle ;
 - r trigger <vsDeliver, src, m >;

rupon event < crash, p > do

- orrect := correct \ { p };
- if flushing = false then
 - flushing := true;
 - trigger <vsBlock>;

- Upon <vsBlockOk> do
 - blocked := true;
 - rtrigger <bebBroadcast, [DSET,view.id,delivered] >;

Upon <bebDeliver, src, [DSET,vid,del] > do

 \checkmark dset:= dset \cup (src,del);

if forall p ∈ correct, (p,mset) ∈ dset then
 trigger <ucPropose, view.id+1, correct, dset >;

- Upon <ucDecided, id, memb, vsdset > do
 - *r* **forall** (p,mset) \in vsdset: p \in memb do
 - *r* **forall** (src,m) ∈ mset: m ∉ delivered **do**
 - \checkmark delivered := delivered $\cup \{m\}$
 - r trigger <vsDeliver, src, m>;
 - view := (id, memb); flushing := blocked :=
 false; dset := delivered := Ø;
 - r trigger <vsView, view>;

Uniform View Synchrony

We now combine the properties of

group membership (Memb1-Memb4) – which is already uniform

uniform reliable broadcast (RB1-RB4) -

VS: A message is vsDelivered in the view where it is vsBroadcast – which is already uniform

Uniform View Synchrony

Using uniform reliable broadcast instead of best effort broadcast in the previous algorithms does not ensure the uniformity of the message delivery



Algorithm 3 (uvsc)

- r upon event < Init > do
 - r view := (0,S);
 - r correct := S;
 - flushing := blocked := false;
 - \checkmark udelivered := delivered := dset := \varnothing ;
 - \checkmark for all m: ack(m) := \emptyset ;

- upon event <vsBroadcast,m) and (blocked = false)
 do
 </pre>
 - \checkmark delivered := delivered $\cup \{m\};$
 - rigger <bebBroadcast,[Data,view.id,m] >;

- upon event < bebDeliver,src,[Data,vid,m]) do
 f (view.id = vid) then
 </pre>
 - \checkmark ack(m) := ack(m) \cup {src};
 - \checkmark if m \notin delivered then
 - \checkmark delivered := delivered $\cup \{ m \}$
 - r trigger <bebBroadcast, [Data,view.id,m] >;

- ✓ upon event (view ≤ ack(m)) and (m ∉ udelivered)
 do
 - udelivered := udelivered ∪ { m }
 - rtrigger <vsDeliver, src(m), m >;

rupon event < crash, p > do

- orrect := correct \ { p };
- if flushing = false then
 - flushing := true;
 - trigger <vsBlock>;

Upon <vsBlockOk> do

- blocked := true;
- frigger <bebBroadcast,
 [DSET,view.id,delivered] >;
- Upon <bebDeliver, src, [DSET,vid,del] > do
 - r dset:= dset \cup (src,del);
 - ✓ if forall p ∈ correct, (p,mset) ∈ dset
 then trigger <ucPropose, view.id+1,
 correct, dset >;

- Upon <ucDecided, id, memb, vsdset > do
 - *r* **forall** (p,mset) \in vs-dset: p \in memb do
 - ✓ forall (src,m) ∈ mset: m ∉ udelivered do
 - \checkmark udelivered := udelivered $\cup \{m\}$
 - r trigger <vsDeliver, src, m>;
 - view := (id, memb); flushing := blocked :=
 false; dset := delivered := udelivered := Ø;
 - r trigger <vsView, view>;